

# **Productions Application Fraud Analysis**

**Prepared by:**

**Team 4**

Wenzhen Zhao, Ying Liu, Ting Gu, Xuan Zhang, Po-Nien Chiang, Yue Shi, Xingjian  
Zheng

**Instructor:** Professor Coggeshall

**Fraud Analytics (DSO 562)**

# Table of Contents

**Executive Summary**----- 2

**Part I: Data Overview**----- 3

**Part II: Variable Construction**----- 10

**Part III: Feature Selection**----- 13

**Part IV: Fraud Algorithm** ----- 15

**Part V: Results** ----- 18

**Part VI: Conclusion** ----- 19

**Appendix** ----- 20

## ❖ Executive Summary

The report provides fraud detection analysis of Credit Card Application Data using supervised algorithms. Main data processing tools are Python, R and Microsoft Power BI.

The original dataset contains 94866 rows of application records with 9 variables of applicants' personal information. The pipeline of our works includes variable construction, data cleaning, feature selection using KS, fraud algorithm implementation, score calculation and conclusion.

During the fraud algorithm implementation process, we fit 6 algorithms with the training set one by one, including Logistic Regression, Gaussian Naive Bayes, XGBoost, Decision Tree, Random Forest and Support Vector Machine models. Then, fraud score is calculated for each record. By sorting all the records in descending order of the scores and subsetting top 10% records for each model, we got the fraud detection rate for each algorithm and found that Random Forest performs the best among all the algorithms.

Using the best-performing model, the fraud detection rate is 80.59% when checking top 1% of the whole population, and the false positive rate is 24%.

## ❖ Part I: Data Overview

### 1.1 Data Summary

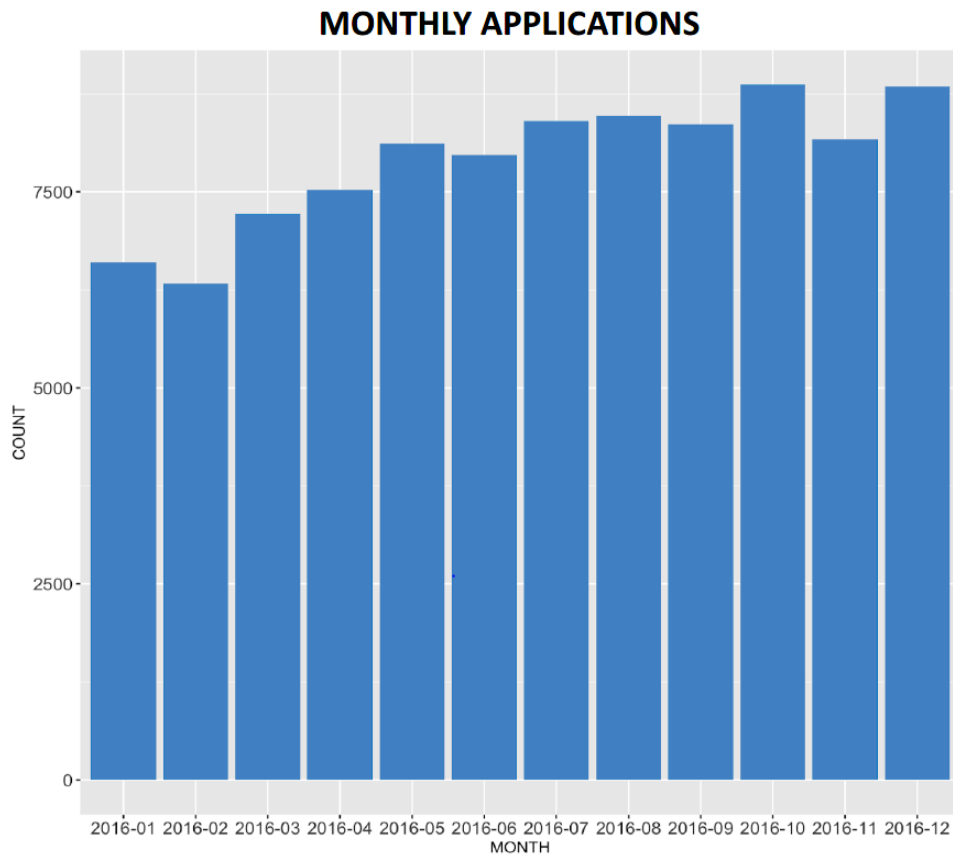
The credit application information in 2016 has 10 variables and 94866 observations in total. Except for the variable “*record*”, all the other variables are categorical variables, including application date, SSN, first name, last name, address, zip code, date of birth, home phone and fraud label.

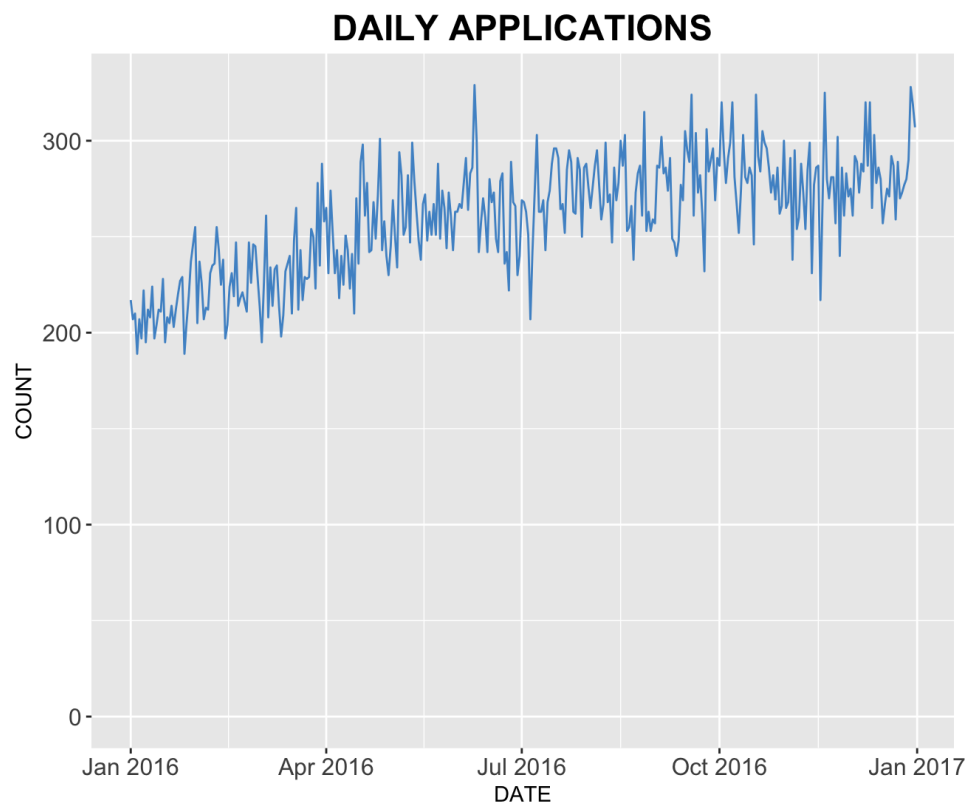
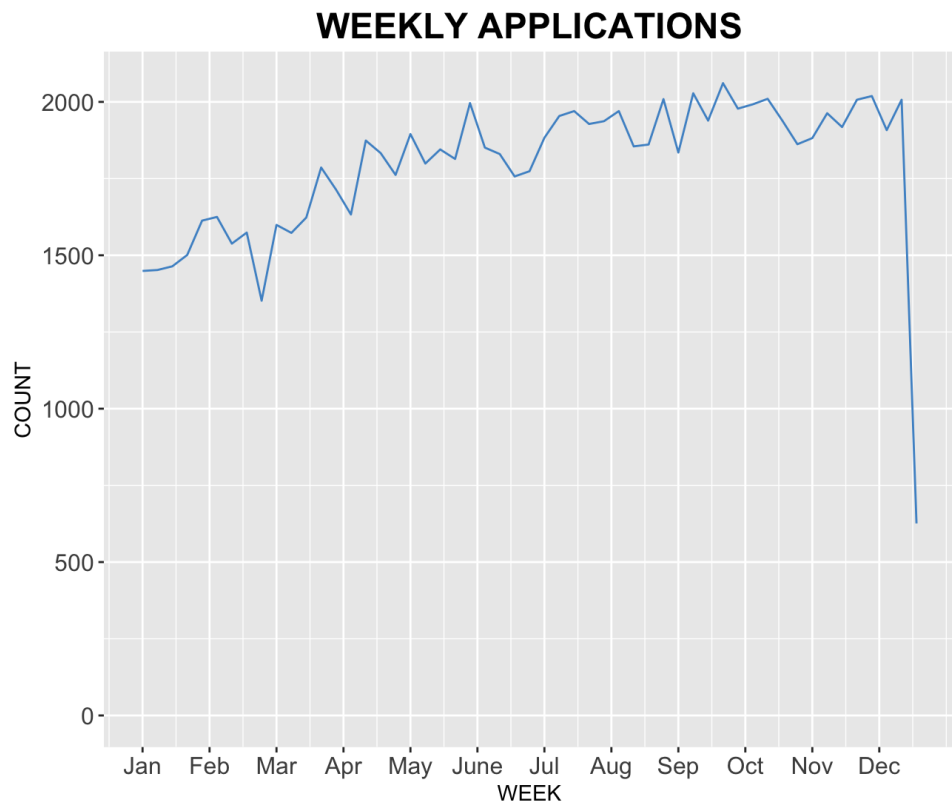
Below is an excerpt of the variables. The complete Data Quality Report is attached in the appendix.

### 1.2 Important variables

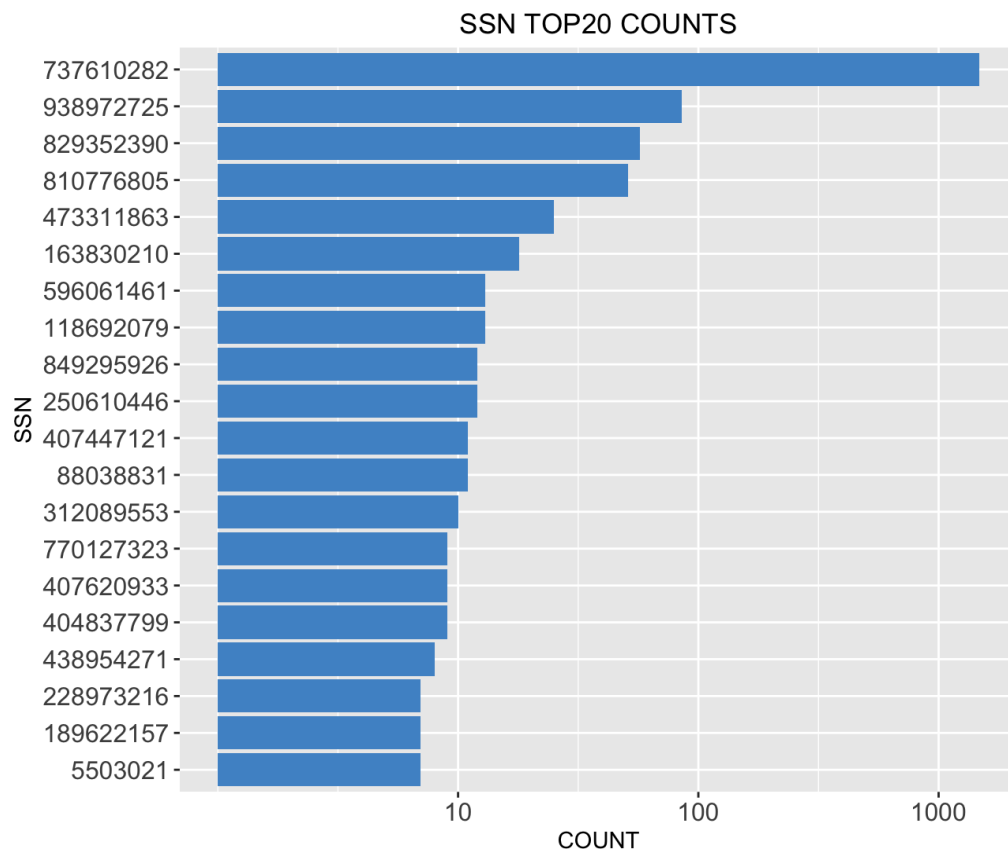
#### 1.2.1 Base variables

**DATE**(categorical variable): The date of applications. There are 365 unique values from 01/01/2016 to 12/31/2016. No missing values exist. Monthly, weekly and daily distributions are shown below.

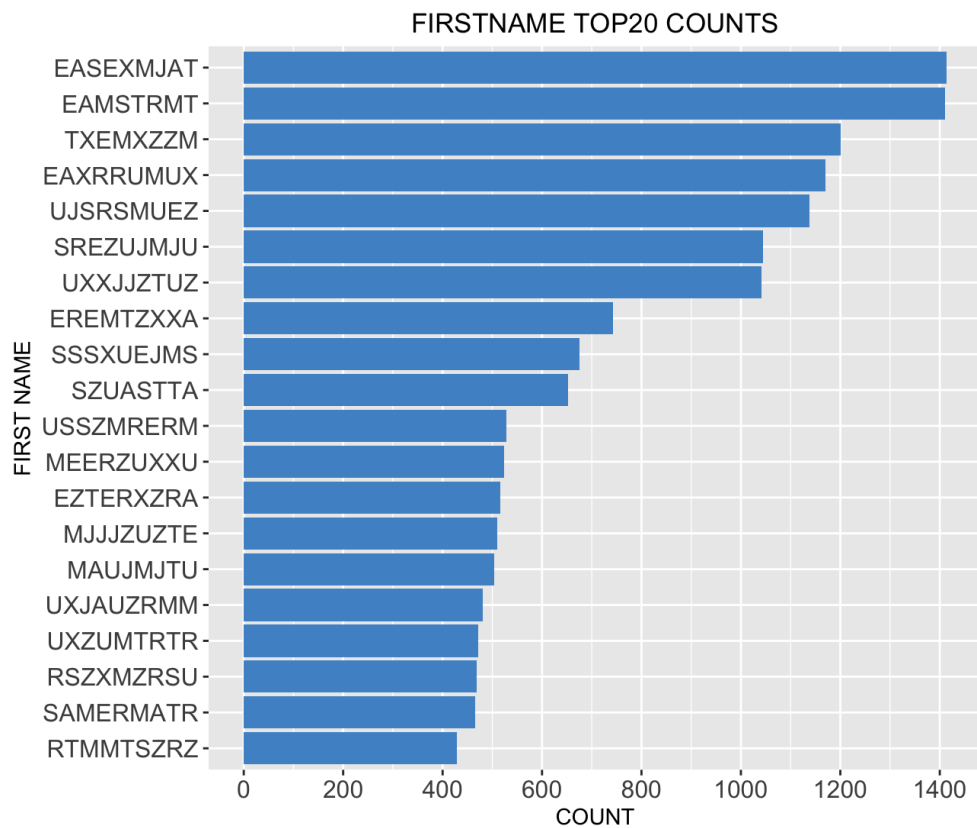




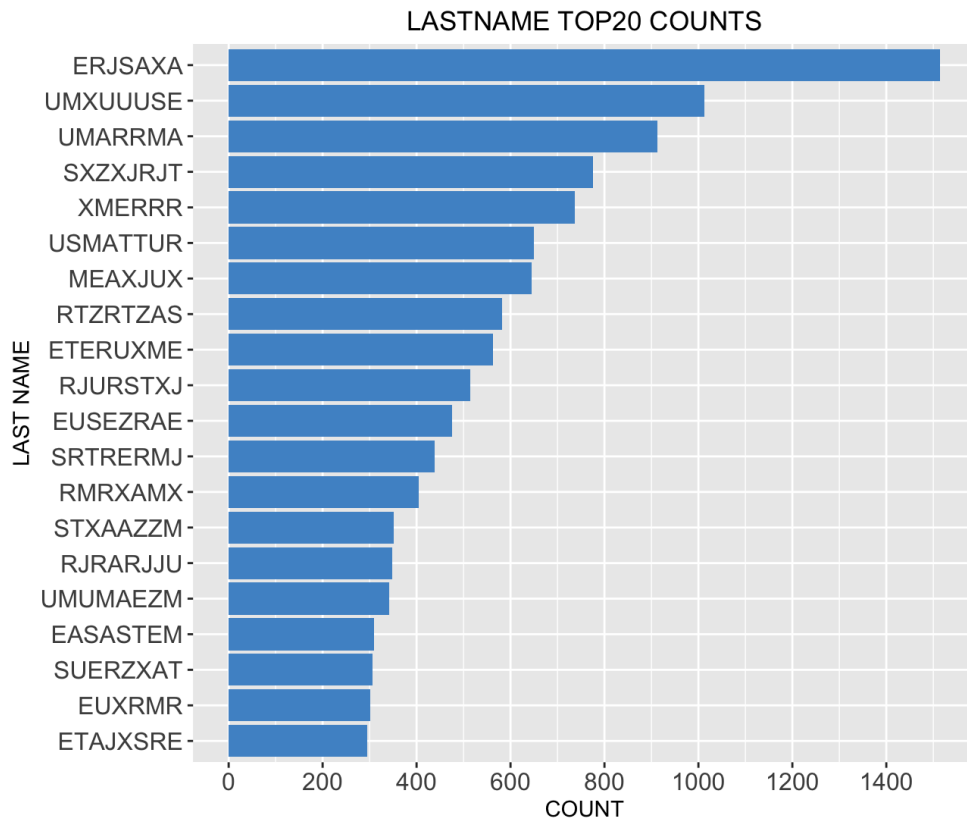
**SSN** (categorical variable): Social Security Number of applicants. There are 86771 unique values. No missing values exist.



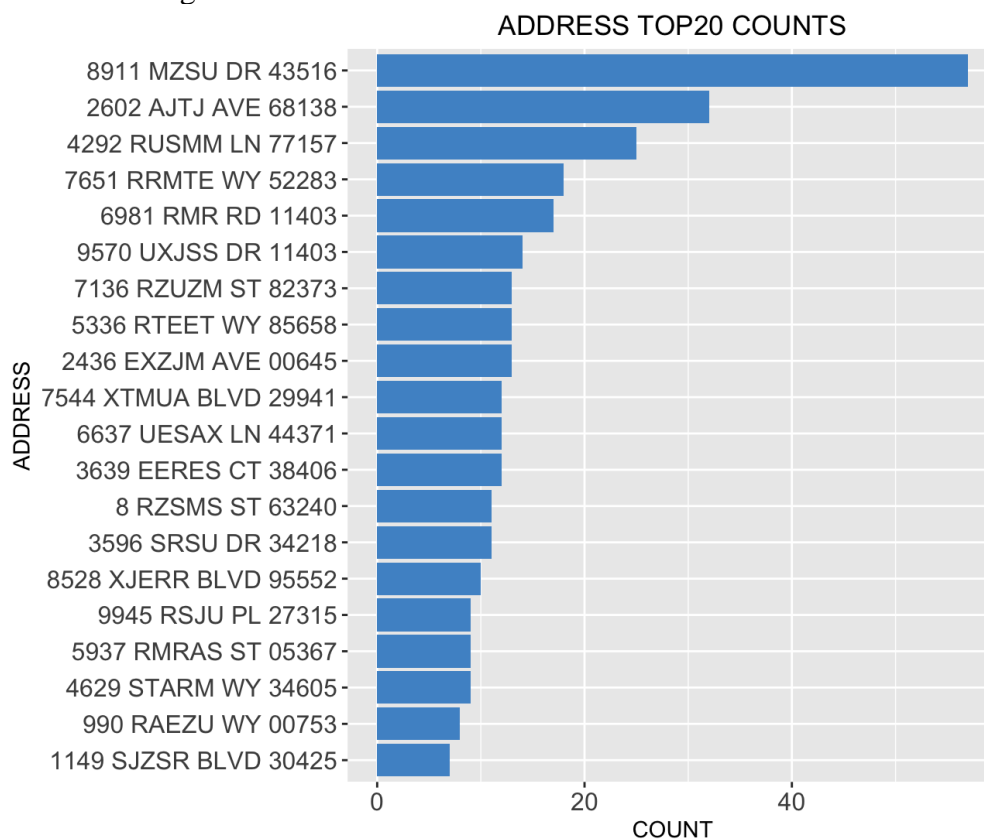
**FIRSTNAME** (categorical variable): First name of applicants. There are 14626 unique values. No missing values exist.



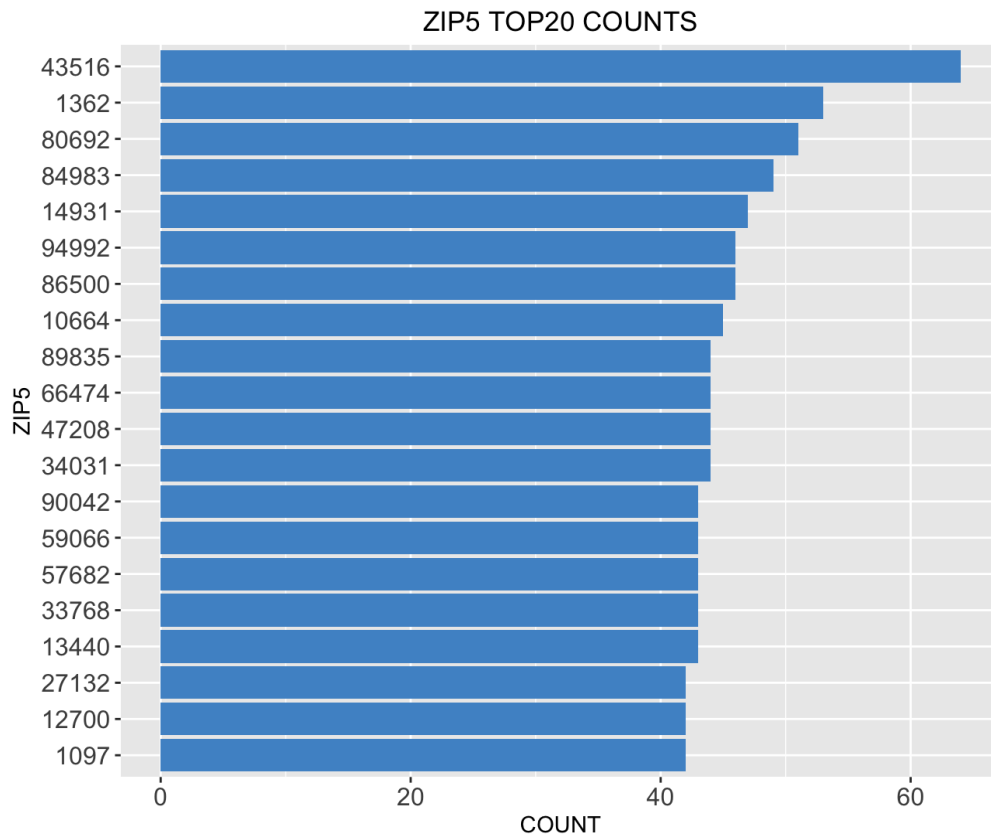
**LASTNAME** (categorical variable): Last name of applicants. There are 31513 unique values. No missing values exist.



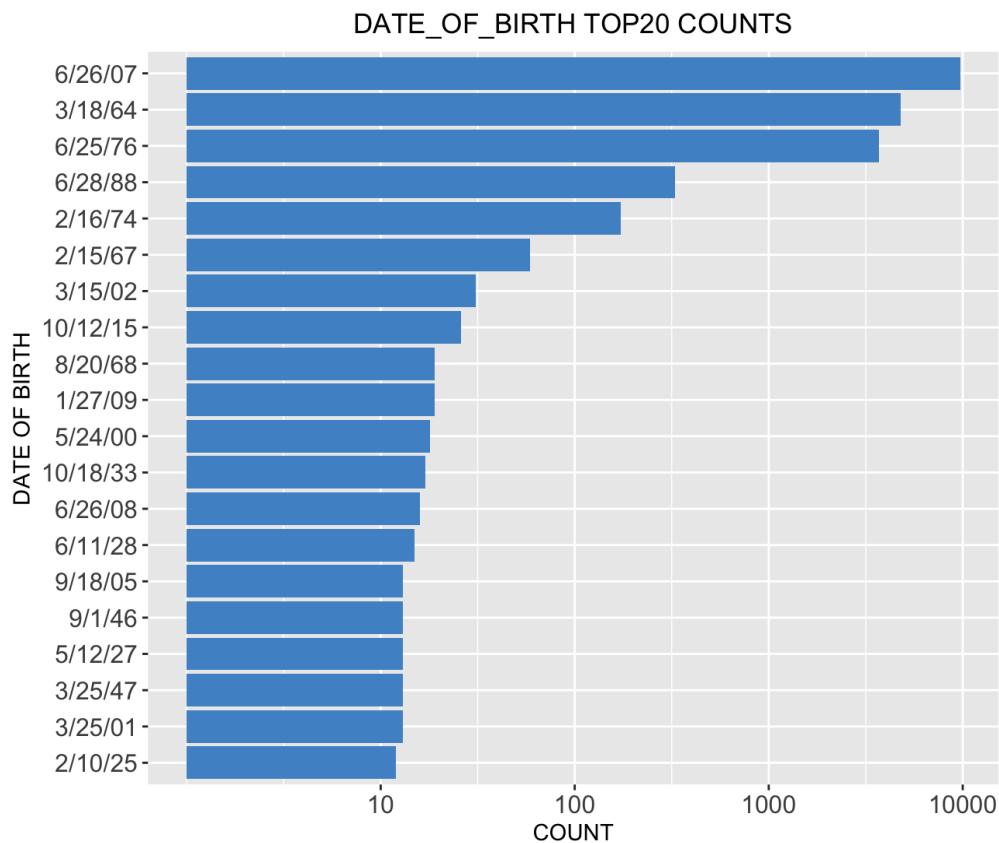
**ADDRESS** (Categorical variable): Address of applicants. There are 88167 unique values and no missing value.



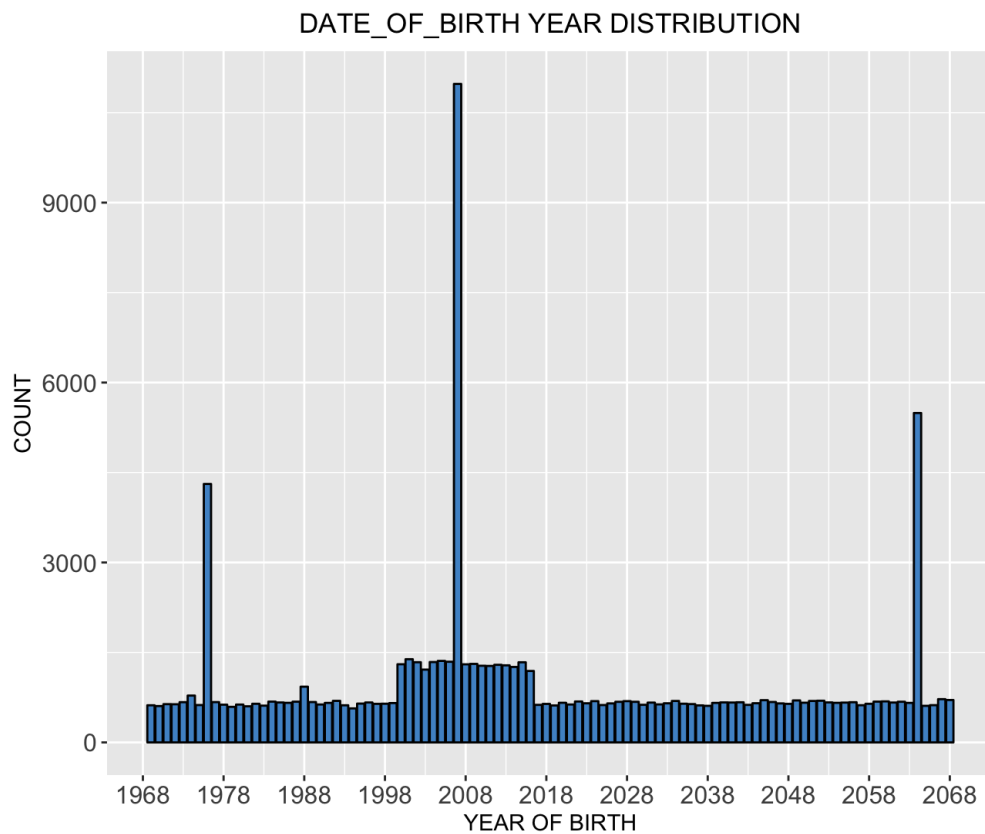
**ZIP5** (Categorical variable): The zip code of applicants' location. There are 15855 unique values and no missing value.



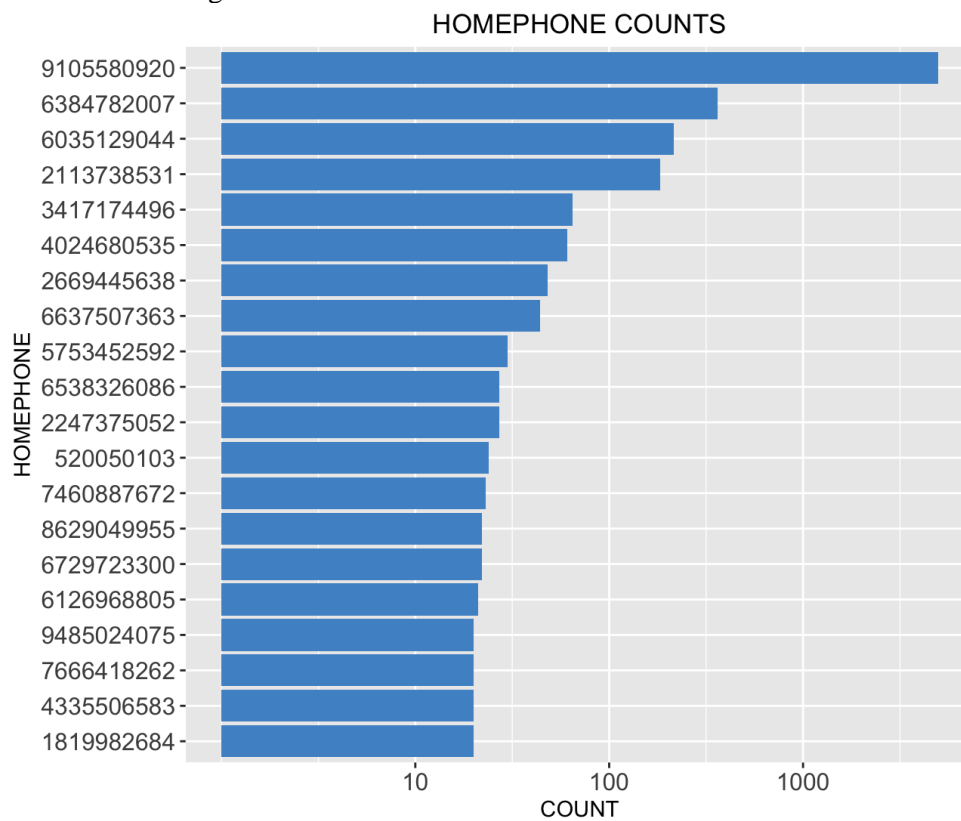
**DOB** (Categorical variable): The date of birth of applicants. There are 30699 unique values and no missing value. The histogram of the distribution of the year when applicants were born is also shown.





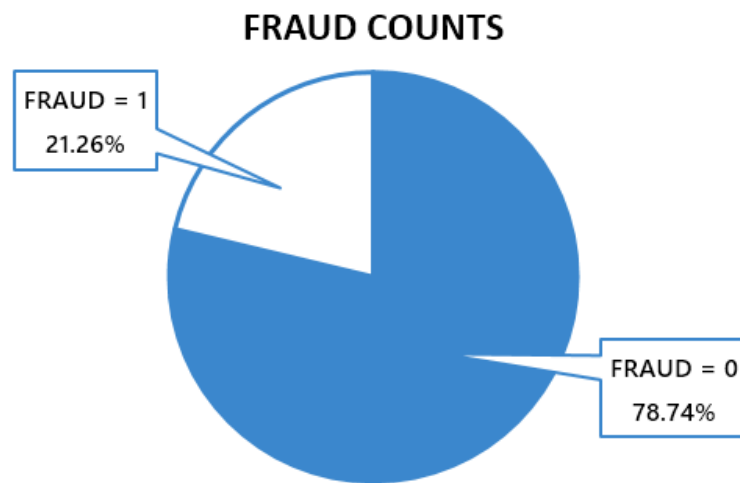


**HOMEPHONE** (Categorical variable): The date of birth of applicants. There are 30699 unique values and no missing value.



### 1.2.2 Fraud labels

**FRAUD** (Categorical variable): Label of fraud definition. There are 2 unique values. Label “1” means fraudulent and label “0” means believable. No missing values exist.



## ❖ Part II: Variable Construction

To quantify the activities and connect these activities to frauds, we derived two types, 20 variables in total from the base variables. Next, by taking chronological order into consideration, we built 5 time windows for each variable and count the appearance of the same values from the previous records within the given time window. The time windows are: **1 day, 3 days, 7 days, 15 days, and 30 days**.

In total, we will have **100** numeric variables from the variable construction section. After initial construction, we reset frivolous values into the baseline counts to avoid including noisy information to our variables.

### 2.1 Type I Variable

In this section, we keep base variables as their original form, and derive the 5 time window counts from these variables.

**ssn\_n**: the number of specific applicant under one particular **ssn** in time windows n.

**homephone\_n**: number of specific applicant under one particular **homephone** in time windows n.

**zip5\_n**: number of specific applicant under one particular **zip\_5** in time windows n.

**dob\_n**: number of specific applicant under one particular **dob** in time windows n.

**address\_n**: number of specific applicant under one particular **address** in time windows n.

By deriving 1 day, 3 days, 7 days, 15 days and 30 days counts from the base variables above, we will construct **25 Type I** variables.

### 2.2 Type II Variable

In this section, we combine original variables to generate three variables bases: two-variable base, three-variable base, and four-variable base. That is, we combine different original variables together as an indicator of different individual applicant, and then calculate the number of this individual in different time windows. In total, there are **75 Type II** variables.

#### 2.2.1 Two-variable Base Variable

**firstname\_lastname\_n**: the number of specific applicant has particular **firstname** and **lastname** in time windows n.

**homephone\_address\_n**: the number of specific applicant has particular **homephone** and **address** in time windows n.

**homephone\_zip5\_n**: the number of specific applicant has particular **Homephone** and **zip5** in time windows n.

**address\_zip5\_n**: the number of specific applicant has particular **address** and **zip5** in time windows n.

**ssn\_address\_n**: the number of specific applicant has particular **ssn** and **address** in time windows n.

**ssn\_dob\_n**: the number of specific applicant has particular **ssn** and **dob** in time windows n.

**ssn\_zip5\_n**: the number of specific applicant has particular **ssn** and **zip5** in time windows n.

**ssn\_homephone\_n**: the number of specific applicant has particular **ssn** and **homephone** in time windows n.

### 2.2.2 Three-variable Base Variable

**firstname\_lastname\_dob\_n**: the number of specific applicant has particular **firstname**, **lastname**, and **dob** in time windows n.

**firstname\_lastname\_ssn\_n**: the number of specific applicant has particular **firstname**, **lastname**, and **ssn** in time windows n.

**firstname\_lastname\_homephone\_n**: the number of specific applicant has particular **firstname**, **lastname**, and **homephone** in time windows n.

**firstname\_lastname\_address\_n**: the number of specific applicant has particular **firstname**, **lastname**, and **address** in time windows n.

**firstname\_lastname\_zip5\_n**: the number of specific applicant has particular **firstname**, **lastname**, and **zip5** in time windows n.

**address\_zip5\_homephone\_n**: the number of specific applicant has particular **address**, **zip5**, and **homephone** in time windows n.

### 2.2.3 Four-variable Base Variable

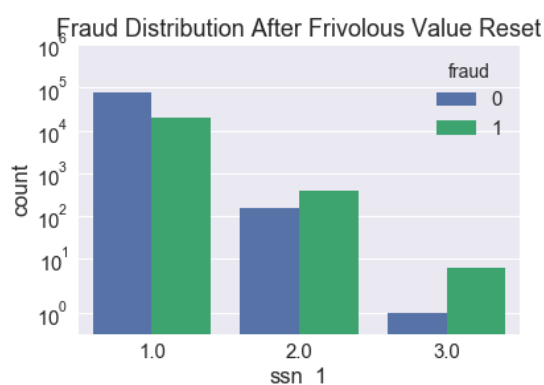
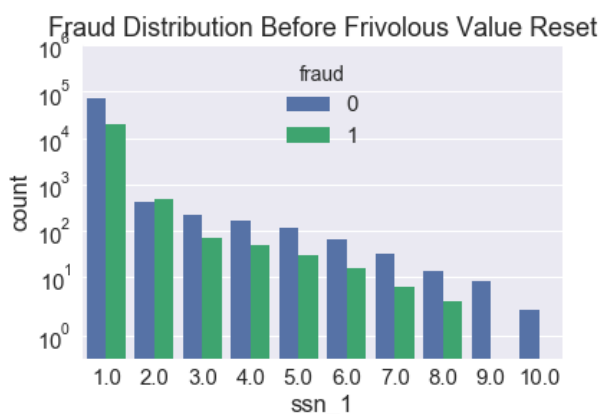
**firstname\_lastname\_ssn\_dob\_n**: the number of specific applicant has particular **firstname**, **lastname**, **ssn**, and **dob** in time windows n.

## 2.3 Reset Frivolous values

After variable construction, we identified frivolous values for **homephone**, **ssn**, **dob** as the outliers in counts. The rationale for this action comes from the understanding of user behavior. For applicants, there are certain patterns of **homephone**, **ssn** and **dob** to be forged when they don't want to reveal these personal informations. Since we cannot distinguish fraud from non fraud based on high counts in these records, we just reset them to the baseline count for the field. The selected frivolous values are shown below.

	0	1	2	3	4
<b>ssn</b>	737610282	None	None	None	None
<b>dob</b>	6/26/07	3/18/64	6/25/76	6/28/88	2/16/74
<b>homephone</b>	9105580920	6384782007	6035129044	2113738531	None

Also, we can tell the functionality of resetting frivolous from visualizations. On the left, the variable **ssn\_1** with frivolous value counts do not separate fraud records at all; while on the right, as **ssn\_1** counts goes up, the proportion of frauds increases significantly.



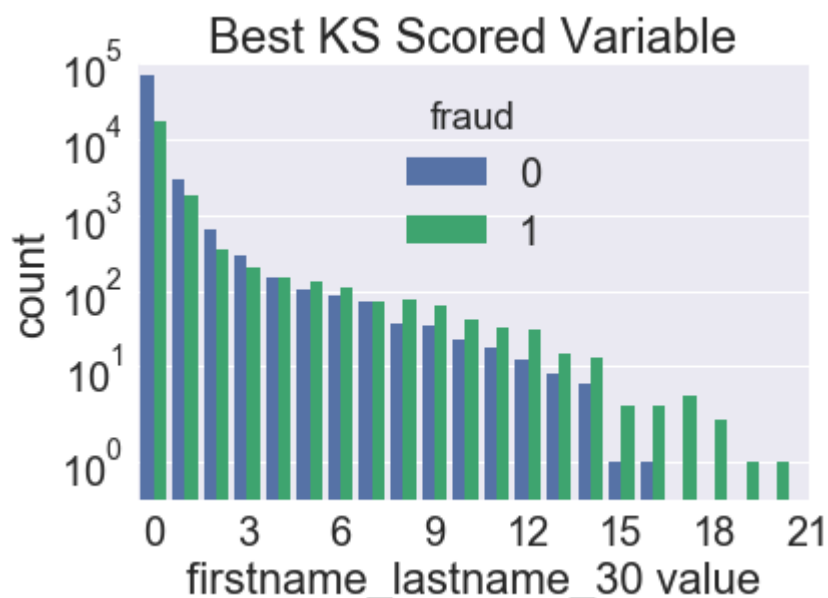
## ❖ Part III: Feature Selection

### 3.1 Variable Selection

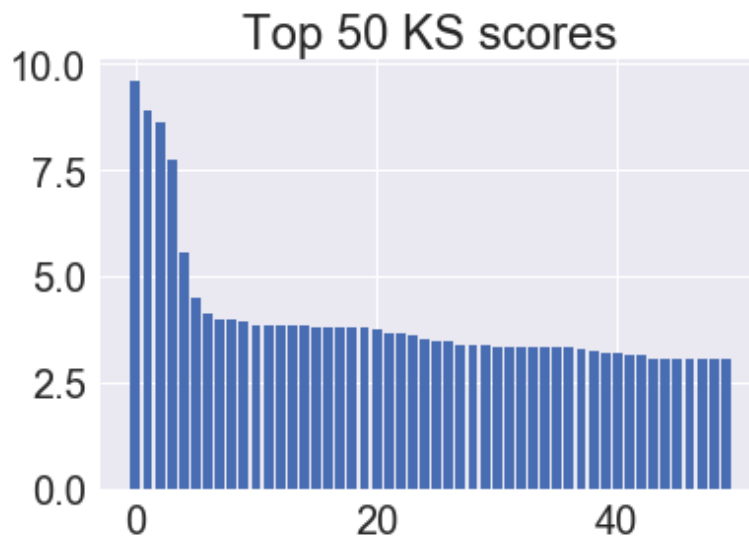
To reduce dimensionality, we performed two-staged feature selection. First we calculated the Kolmogorov Smirnov (KS) scores for each variable and selected the top **50** variables out of **100**. Next we performed forward selection using logistic regression as our baseline model, and selected top **25** variables out of **50** based on recall rates of each combination of variables.

#### 3.1.1 Kolmogorov Smirnov (KS) Score

KS score is widely used as a measurement of the goodness of a variable to separate frauds from non-frauds. To calculate KS score of variables, we derive the cumulative distribution (CDF) of variable values of frauds and non-frauds respectively and take the max difference of CDF, times 100, as the KS score for the variable.



As the plot above shows, '**firstname\_lastname\_30**' is the best KS scored variable out of the 100 variables we created in part 2. The KS score is **9.62**. Below we also visualized the distribution of the top 50 KS scores.



### 3.1.2 Forward selection

After selecting the top 50 KS scored variables, we performed forward selection to further reduce the dimensionality to 25. For interpretability and clarity, we chose logistic regression as our baseline model.

Since the statistics our customer asked for was FDR, which is by definition recall rate ( $\text{true positive} / (\text{true positive} + \text{false negative})$ ), we selected recall rate of each forward stepwise combination of variables as our optimization objective. The top 25 variables selected are shown below.

0	firstname_lastname_1		
1	firstname_lastname_3	13	address_30
2	ssn_homephone_30	14	firstname_lastname_7
3	zip5_15	15	address_zip5_15
4	zip5_30	16	address_zip5_7
5	ssn_homephone_15	17	homephone_zip5_30
6	firstname_lastname_30	18	address_7
7	ssn_7	19	address_zip5_30
8	ssn_15	20	homephone_zip5_15
9	ssn_dob_30	21	address_15
10	ssn_30	22	address_zip5_homephone_30
11	ssn_dob_15	23	firstname_lastname_zip5_7
12	ssn_dob_7	24	ssn_zip5_15

## ❖ Part IV: Fraud Algorithm

### 4.1 Training/Testing/Out of Time Dataset Split

After selecting 25 features, we split the dataset into three parts before training the models: training data, testing data, and out-of-time data. In our case, we took the first 10 months, from January to October, as training and testing data, and then split the training and testing set with the proportion of 8:2 randomly. In addition, the last two months were defined as the out-of-time set. As a result, we got these three sets with similar fraud rate:

Dataset	Number of Records	Number of Frauds	Fraud Rate
Training	62280	12669	20.34%
Testing	15570	3157	20.28%
Out-of-time	17016	4338	25.49%
Total	94866	20164	21.26%

In our models, we used training set to fit the model, and then tune the model according to the testing set. At the end, we tested our model with out-of-time data.

### 4.2 Supervised Models

In total, we tried 6 supervised models. In each model, we used the probability for data belonging to class 1 as fraud score. Then we sorted all the records by scores in descending order. After subsetting top 10% records for each model, we got the fraud detection rate shown as below:

	FDR @10%		
Model	Training	Testing	Out_of_Time
Logistic Regression	18.68%	18.31%	15.35%
Gaussian Naive Bayes	16.35%	16.15%	14.33%
Decision Tree	20.66%	16.57%	10.58%
Random Forest	18.92%	18.21%	16.16%
XGBoost	19.25%	18.49%	15.49%
Support Vector Machine	17.99%	17.99%	14.89%

#### 4.2.1 Logistic Regression



Logistic regression model uses sigmoid function to estimate the probability and predict probability of each record. The fraud detection rates at 10% population for training, testing, and out-of-time dataset are respectively 18.68%, 18.31%, and 15.35%. This indicates that logistic regression works well for this dataset.

#### **4.2.2 Gaussian Naive Bayes**

Gaussian Naive Bayes is a models that assign class labels to problem instances, represented as vectors of feature values, where the class labels are drawn from some finite set. All naive Bayes classifiers assume that the value of a particular feature is independent of the value of any other feature, given the class variable. Under this model, the fraud detection rates at 10% population for training, testing, and out-of-time dataset are 16.35%, 16.15%, and 14.33% respectively. The result doesn't endorse the mode as a good fit.

#### **4.2.3 Decision Tree**

Decision tree model is based on decision tree structure. The tree can be learned by splitting the source set into subsets based on an attribute value test. The attribute value test we used here is "information entropy". The fraud detection rates at 10% population for training, testing and out-of-time dataset are respectively 20.66%, 16.57% and 10.58%. The result demonstrates that this model has some relatively huge gaps. That is, the decision tree model is not a good model for this dataset.

#### **4.2.4 Random Forest**

Random Forest is a modified version of decision tree. It combines multiple small trees, and considers a subset of features to fit the model. Therefore, it performs far better than decision tree on the overfitting problem. As shown in the table above, random forest performs much better on out-of-date set fraud detection rate than decision tree.

Using random forest model, we captured more than 18% frauds, when checking only 10% of the whole population on training and testing sets. Furthermore, the rate for out-of-time is 16.16%.

#### **4.2.5 XGBoost**

XGBoost is abbreviation for "Extreme Gradient Boosting", a machine learning technique for regression and classification problems, which produces a prediction model in the form of an ensemble of weak prediction models, typically decision trees. When XGBoost applied in our analysis, the fraud detection rates at 10% population for training, testing and out-of-time are 19.2%, 18.5% and 15.5 respectively.

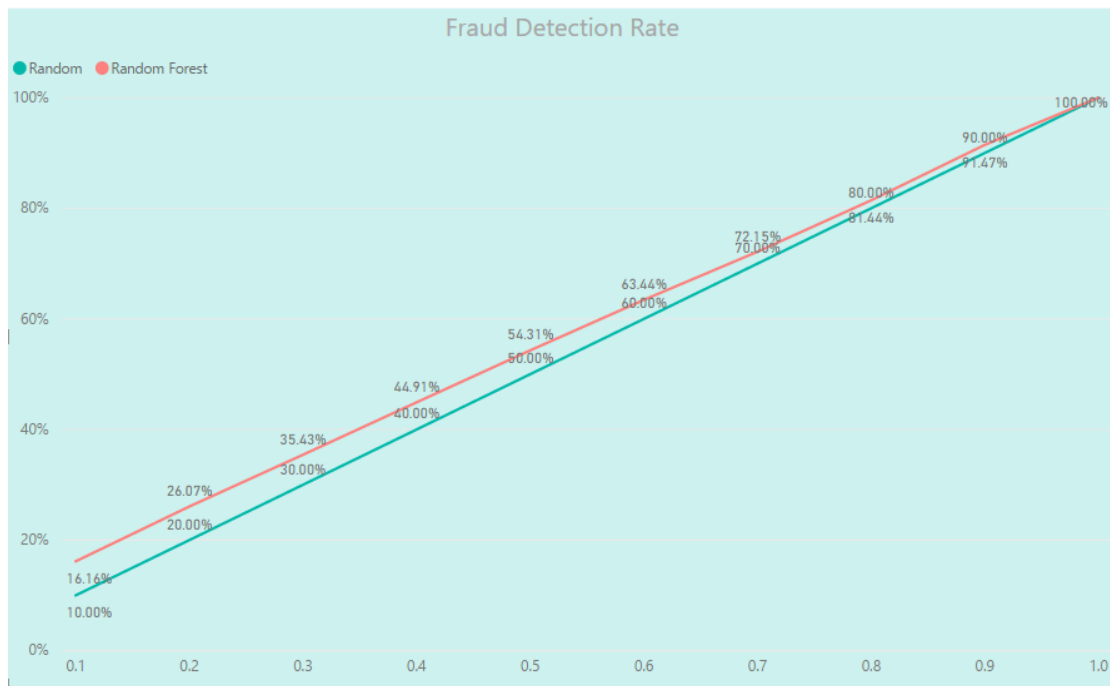
#### **4.2.6 SVM**

SVM is abbreviation for “Support Vector Machines”, a supervised learning models used for classification and regression analysis. SVM training algorithm usually works as a non-probabilistic binary linear classifier, although methods such as Platt scaling exist to use SVM in a probabilistic classification setting, which is what we use to get fraud score.

Under SVM model, we captured 17.99% of the frauds when checking population with top 10% fraud score on training and testing dataset. In the same way, the rate for out-of-time was 14.89%.

## ❖ Part V: Results

According to the analysis above, we concluded that Random Forest performs best among all models. The following graph illustrates the comparison between fraud detection rates generated by Random Forest and random selection. The line of Random Forest is always higher than the line of random selection. This shows Random Forest captured fraud more accurately than random selection.



The following table summarizes both bin statistics and cumulative statistics, when random forest is applied to detect fraud. It contains the information about the number of goods, bads, cumulative goods, cumulative bads, percentage of goods, bads, cumulative goods, cumulative bads, false positive ratio and KS. As shown in table, given the percentage of records we examine, the less suspicious customers we choose, the less number of bads (fraud) can be detected from the model.

Overall Bad Rate: 25.49%		Bin Statistics				Cumulative Statistics				
Population Bin	Total # records	# Good	#Bad	%Good	%Bad	Cumulative Good	Cumulative Bad	%Good	%Bad(FDR)	False Pos Ratio
1%	170	33	137	19.41%	80.59%	33	137	0.26%	3.16%	0.24
2%	170	64	106	37.65%	62.35%	97	243	0.77%	5.60%	0.60
3%	170	107	63	62.94%	37.06%	204	306	1.61%	7.05%	1.70
4%	170	95	75	55.88%	44.12%	299	381	2.36%	8.78%	1.27
5%	170	105	65	61.76%	38.24%	404	446	3.19%	10.28%	1.62
6%	170	98	72	57.65%	42.35%	502	518	3.96%	11.94%	1.36
7%	171	126	45	73.68%	26.32%	628	563	4.95%	12.98%	2.80
8%	170	123	47	72.35%	27.65%	751	610	5.92%	14.06%	2.62
9%	170	124	46	72.94%	27.06%	875	656	6.90%	15.12%	2.70
10%	170	125	45	73.53%	26.47%	1000	701	7.89%	16.16%	2.78
11%	170	122	48	71.76%	28.24%	1122	749	8.85%	17.27%	2.54
12%	170	121	49	71.18%	28.82%	1243	798	9.80%	18.40%	2.47
13%	171	132	39	77.19%	22.81%	1375	837	10.85%	19.29%	3.38
14%	170	126	44	74.12%	25.88%	1501	881	11.84%	20.31%	2.86
15%	170	133	37	78.24%	21.76%	1634	918	12.89%	21.16%	3.59
16%	170	130	40	76.47%	23.53%	1764	958	13.91%	22.08%	3.25
17%	170	123	47	72.35%	27.65%	1887	1005	14.88%	23.17%	2.62
18%	170	131	39	77.06%	22.94%	2018	1044	15.92%	24.07%	3.36
19%	171	131	40	76.61%	23.39%	2149	1084	16.95%	24.99%	3.28
20%	170	123	47	72.35%	27.65%	2272	1131	17.92%	26.07%	2.62
21%	170	126	44	74.12%	25.88%	2398	1175	18.91%	27.09%	2.86
22%	170	128	42	75.29%	24.71%	2526	1217	19.92%	28.05%	3.05

## Part VI: Conclusion

From comparing all the above models and their performances, we demonstrate that Random Forest Model with the 25 variables selected during the forward feature selection process performed best, and presents significantly more accurate result than random selection in the fraud detection process.

## ❖ Appendix

### Data Quality Report

#### 1. Summary Statistics

In the applications dataset, there are 10 variables and 94866 observations in total. Except for the variable “record”, all the other variables are categorical variables. Therefore, the summary data is provided below, including the number of records for each variable, the unique value of each variable and the populated rate.

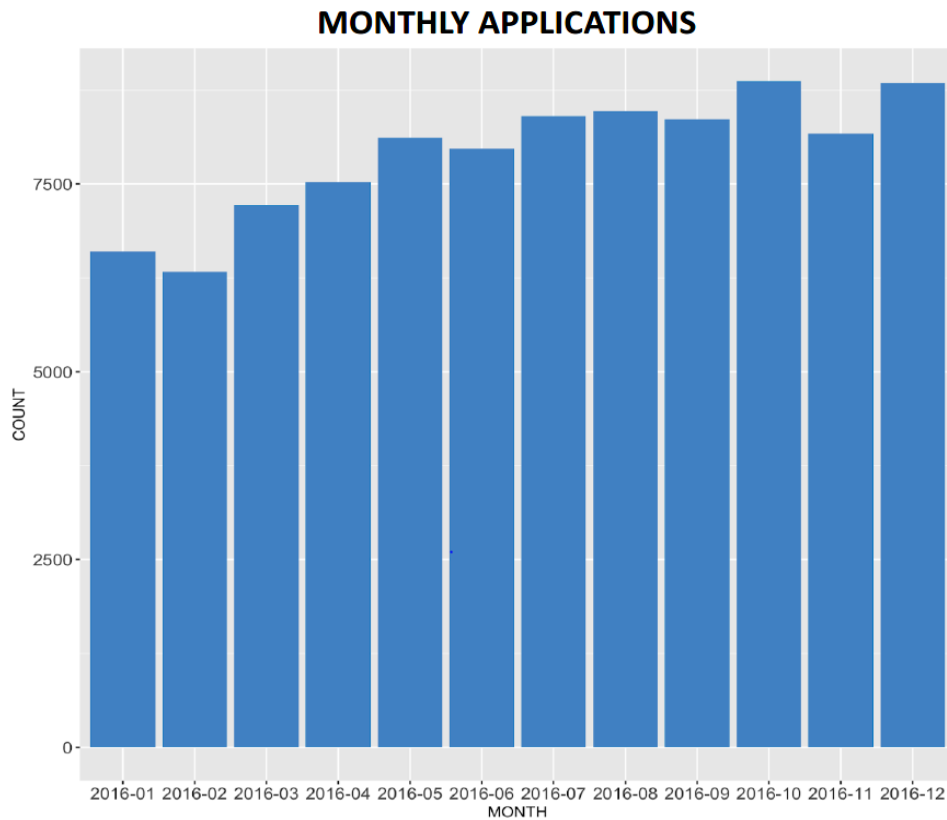
<i>Field Name</i>	<i># of Unique Values</i>	<i>% of Populated</i>
RECORD	94866	100
DATE	365	100
SSN	86771	100
FIRSTNAME	14626	100
LASTNAME	31513	100
ADDRESS	88167	100
ZIP5	15885	100
DOB	30599	100
HOMEPHONE	20762	100
FRAUD	2	100

#### 2. Detailed Information for Each Field

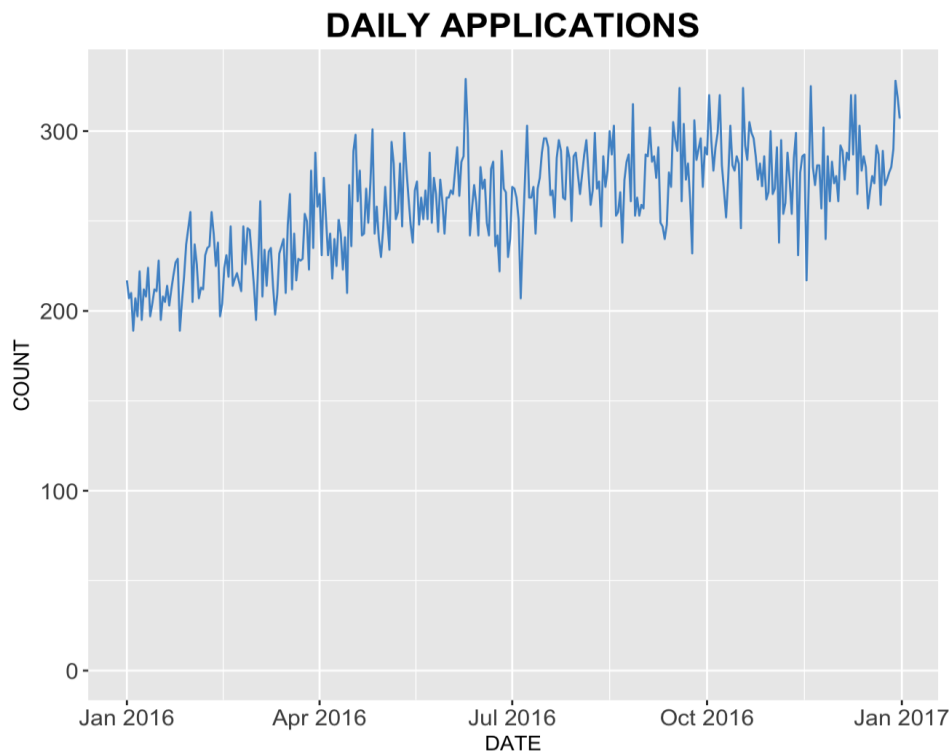
<i>Field Name</i>	<i>Description</i>
<b>RECORD</b>	The number of each record. Discrete data with metric.

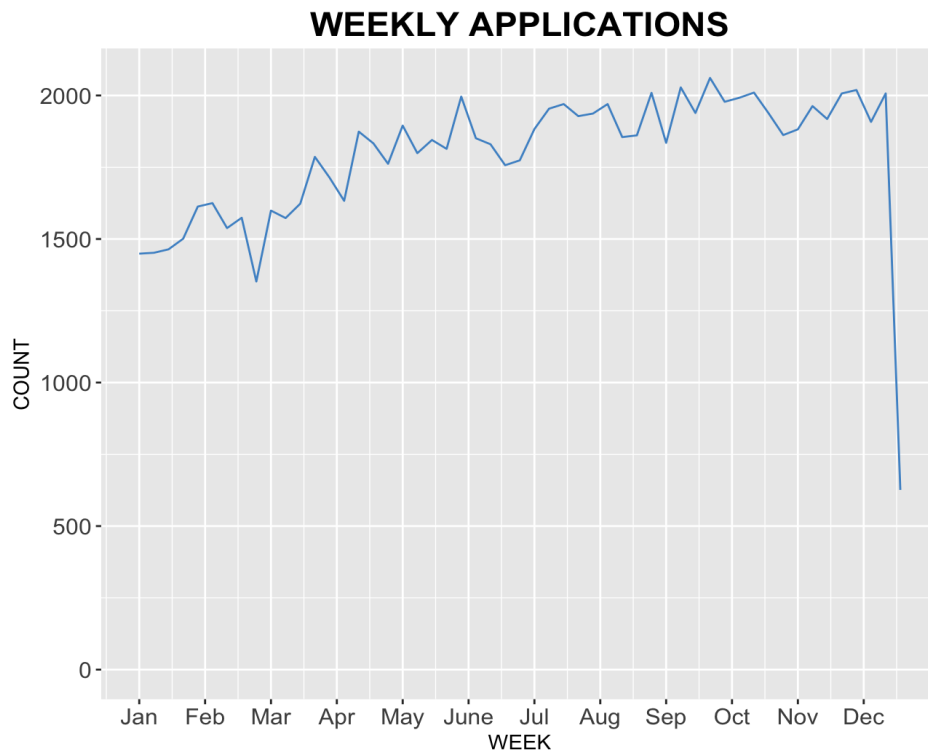
<i>Field Name</i>	<i>Description</i>
<b>DATE</b>	The date of applications. Discrete date data without metric.

In order to better describe the distribution, we visualized the monthly distribution.



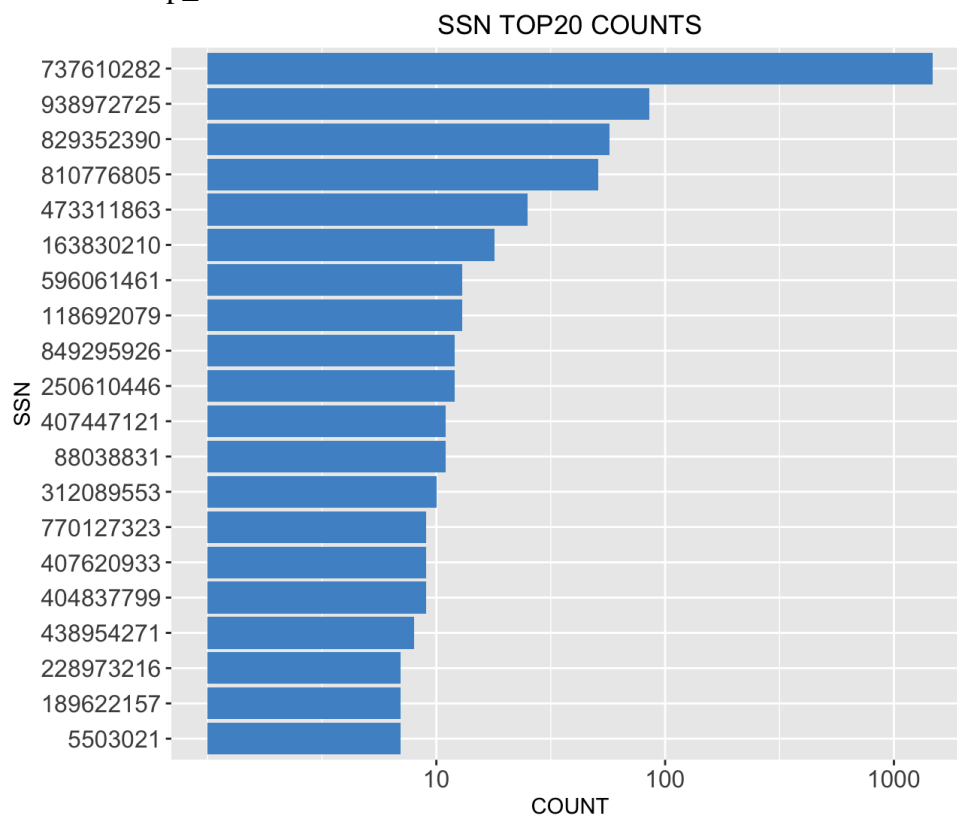
Then we further divided monthly data into weekly data and daily data, and we plotted daily and weekly applications in the following line charts.



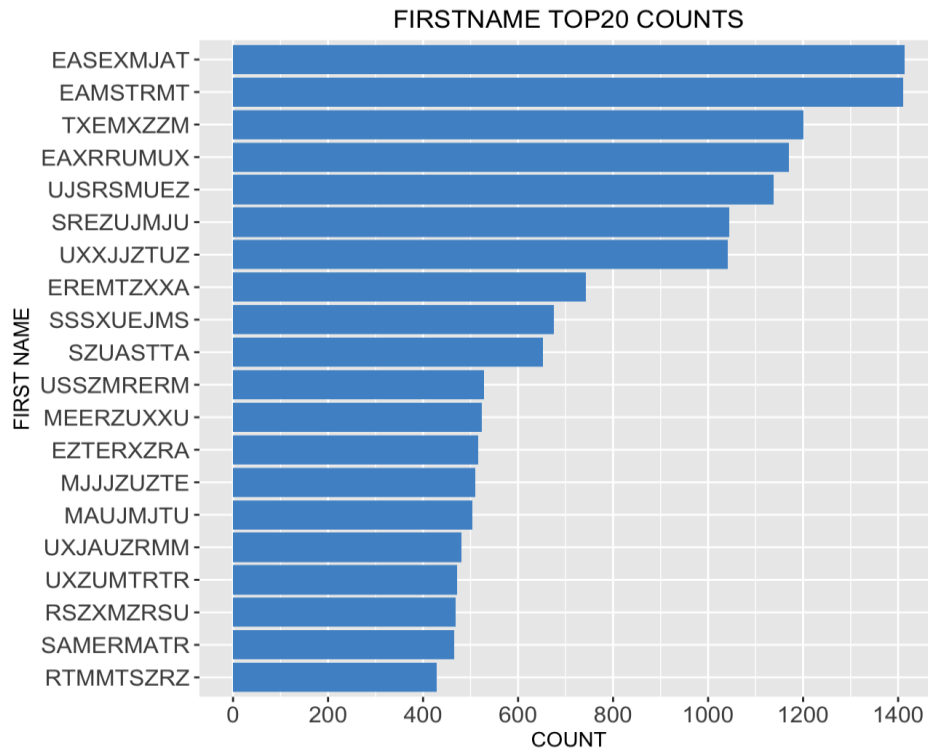


Field Name	Description
SSN	Social Security Number of applicants with 9-digit numerical value. Categorical variable without metric.

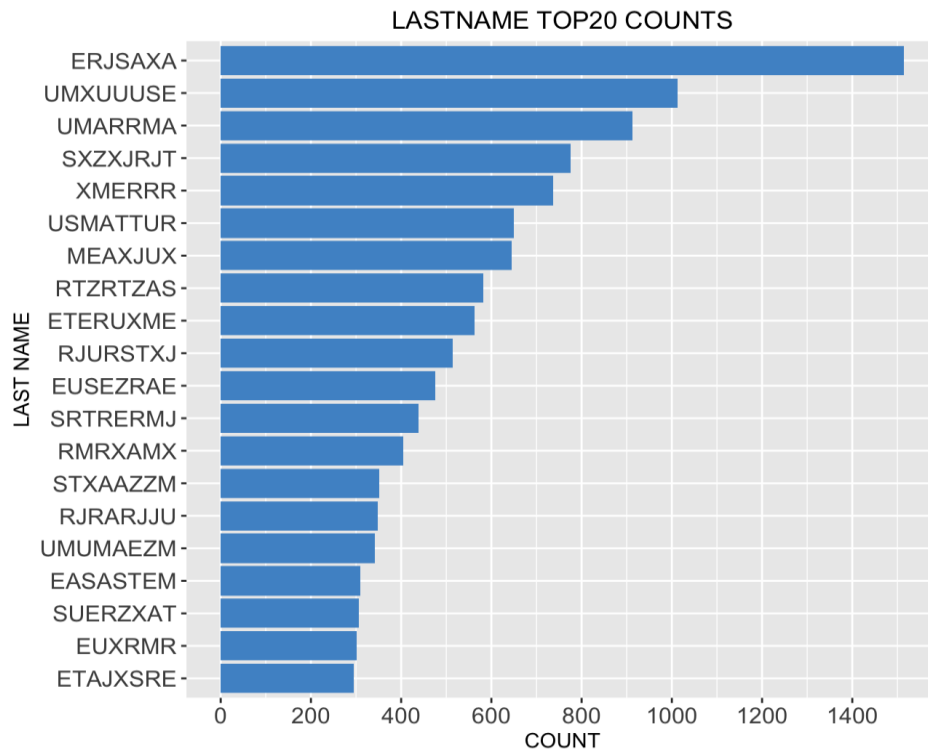
Due to the large difference among the SSN counts, we used log-scale of the count to better visualize top\_20 distribution.



<i>Field Name</i>	<i>Description</i>
<b>FIRSTNAME</b>	First name of applicants. Categorical variable without metric.

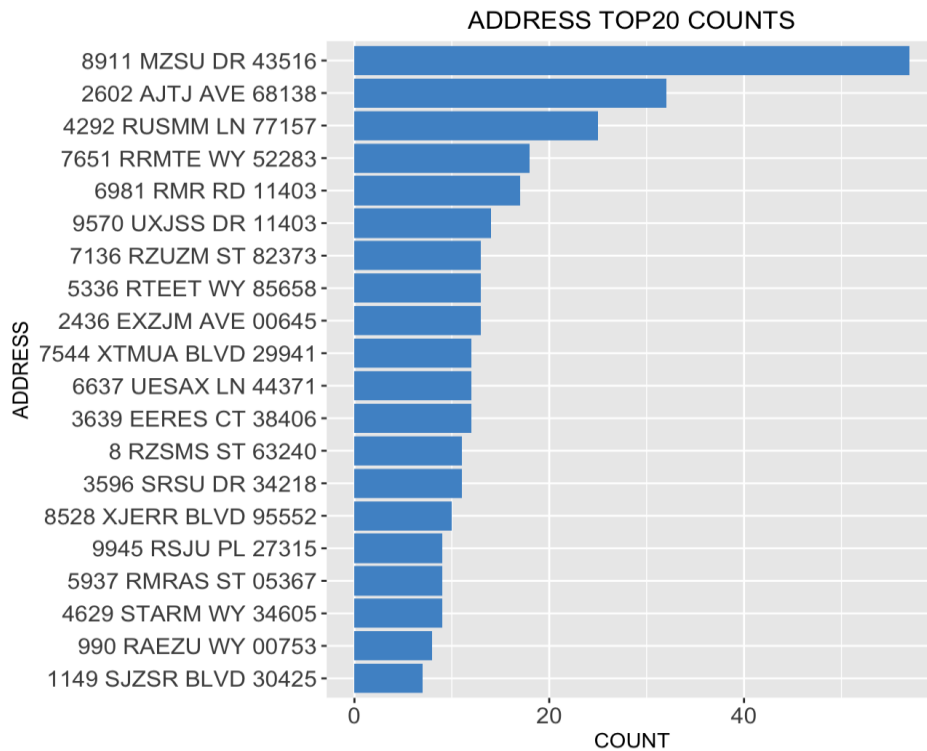


<i>Field Name</i>	<i>Description</i>
<b>LASTNAME</b>	Last name of applicants. Categorical variable without metric.

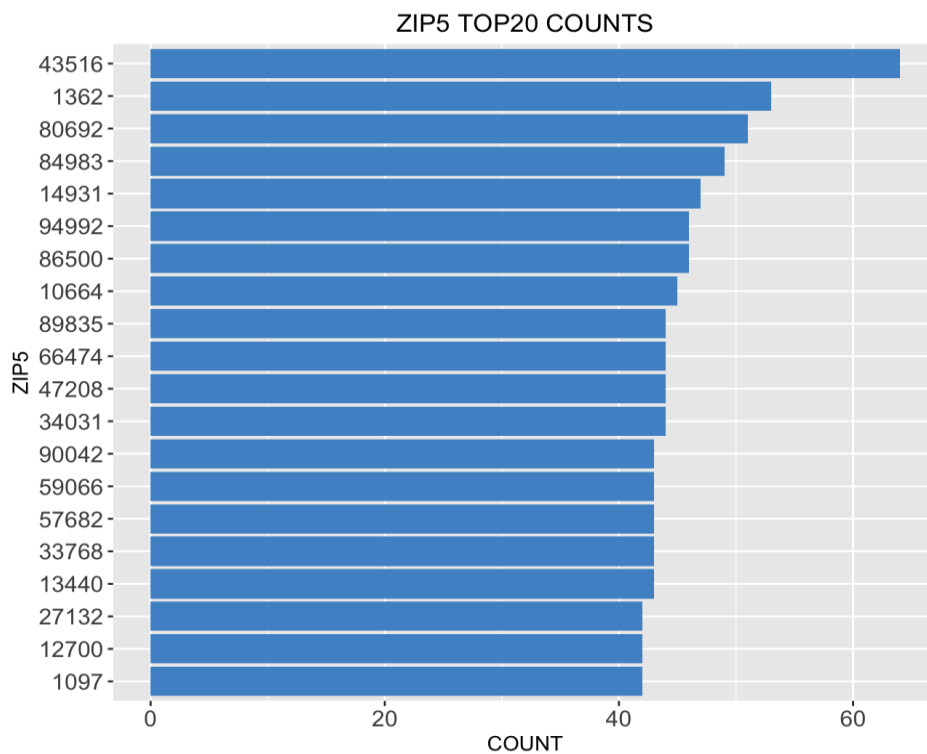




<i>Field Name</i>	<i>Description</i>
<b>ADDRESS</b>	Address of applicants with alphanumeric string. Categorical variable.

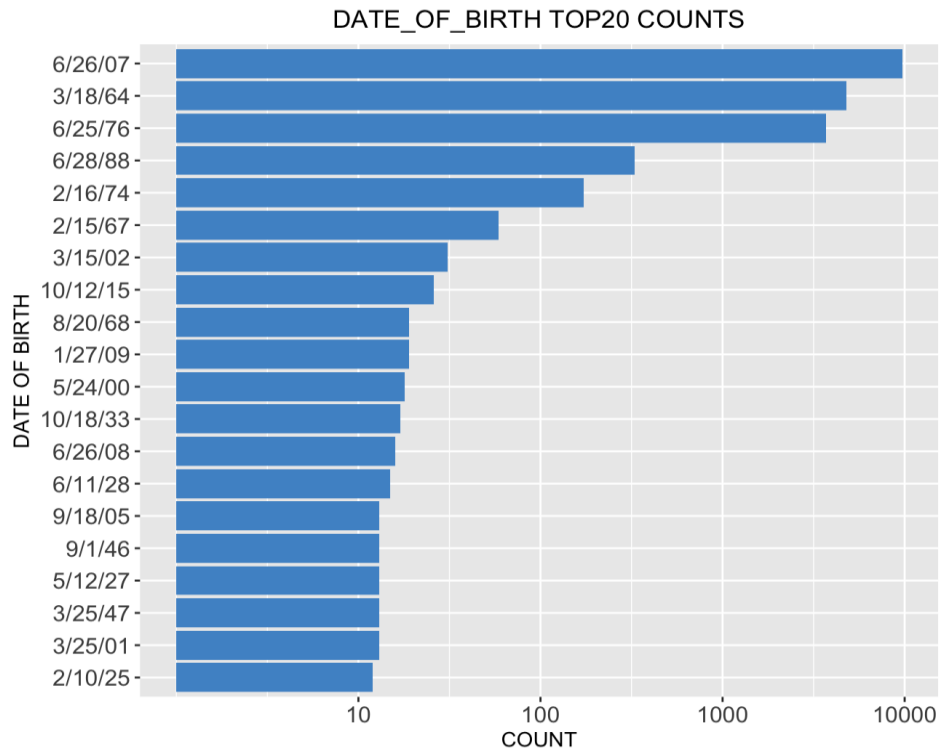


<i>Field Name</i>	<i>Description</i>
<b>ZIP5</b>	The zip code of applicants' location with 5-digit number value. Categorical variable without metric.

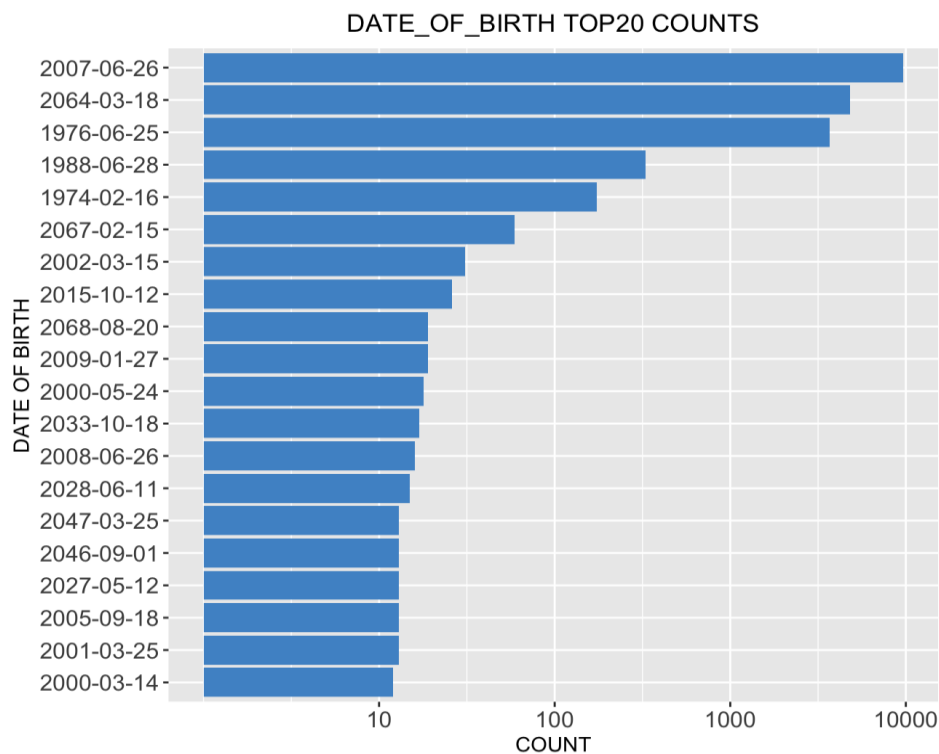


Field Name	Description
<b>DOB</b>	The date of birth of applicants. Categorical variable without metric.

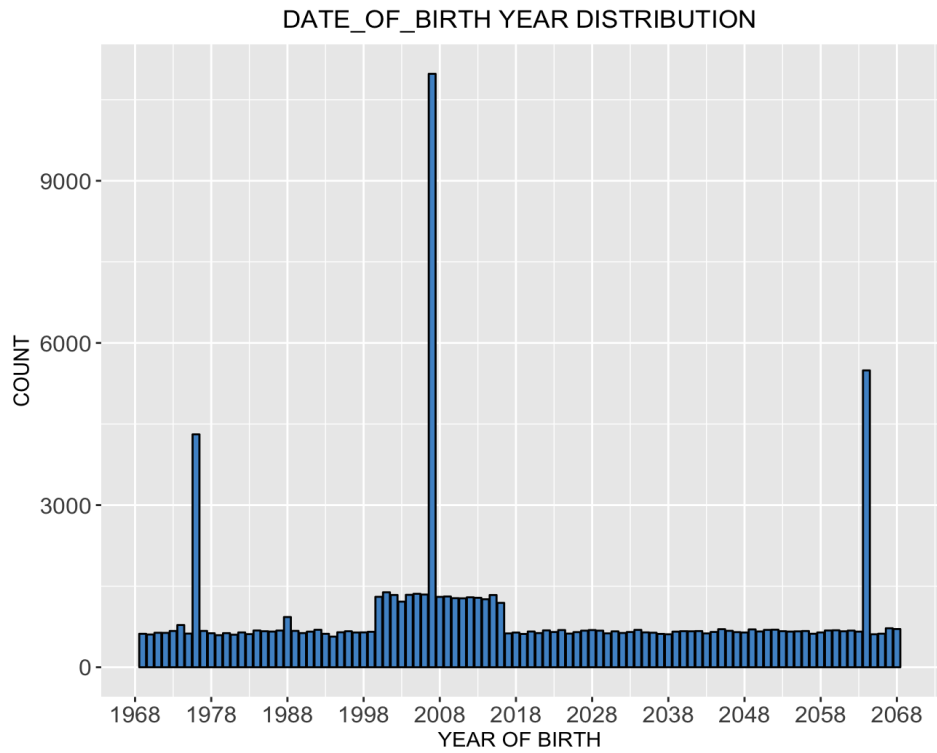
We used log-scale of the count for date of birth to better visualize top\_20 distribution.



Then, we converted the date to year-month-day format to visualize the top\_20 birthdates.

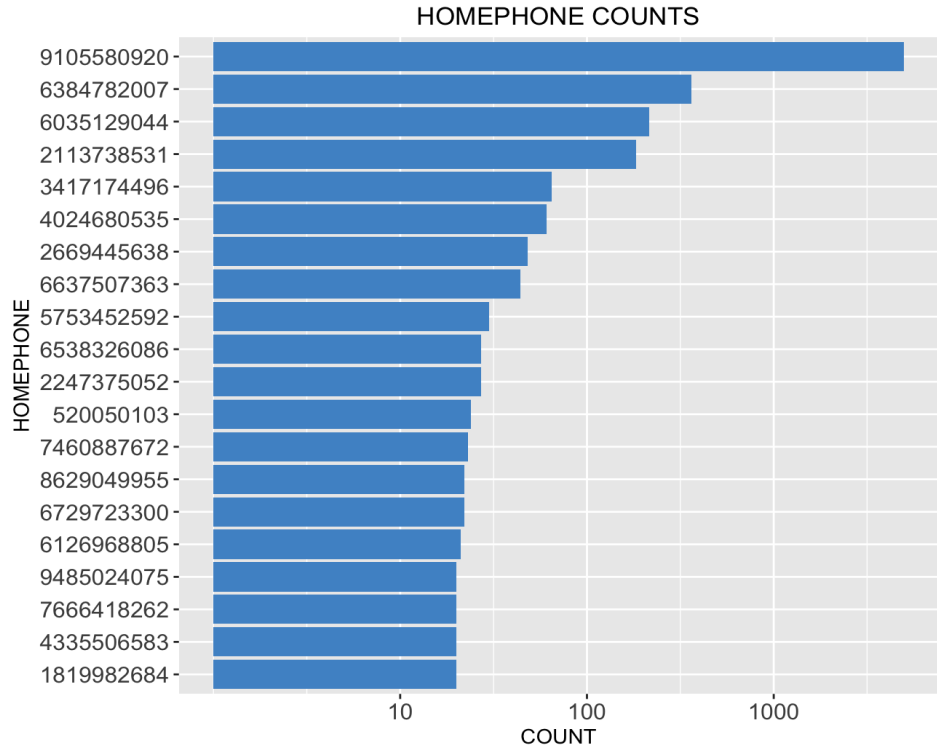


Furthermore, we plotted the histogram of the distribution of the year when applicants were born.



Field Name	Description
HOMEPHONE	Home phone number of applicants. Categorical variable without metric.

We also used log-scale technique to visualize the home phone counts better below.



<i>Field Name</i>	<i>Description</i>
<b>FRAUD</b>	Label of fraud definition. Categorical variable.

